Scene geometry

Disoccluded content inpainting

Motion forecast

Multi Plane Images (MPIs)

Novel view

$I_{t+1}$

Generalized MPI Representation

Entangled motion forecast

Entangled scene representation

Person cut across layers

VEST-MPI

$t$

$I_{t-1}$

$I_t$

Ego-motion forecast

Object motion forecast

Disentangled motion forecast

Estimated depth

Disoccluded content inpainting

Disentangled scene representation in 3D space

Novel view

$I_{t+1}$

Ours

$t$

$I_{t-1}$

$I_t$

Scene geometry

Disoccluded

content inpainting

Motion forecast

Generalized MPI Representation

Multi Plane Images (MPIs)

Novel view

$I_{t+1}$

Person cut
across
layers

VEST-MPI

Entangled scene representation

Quantized scene geometry

Motion
forecast

Estimated
depth

Disoccluded
content inpainting

Novel view

$I_{t+1}$

Continuous scene geometry

Disentangled scene representation in 3D space

Ours

$t$

$I_{t-1}$

$I_t$

Multi Plane Images (MPIs)

Person cut across layers

Generalized MPI Representation

Novel view

Disoccluded content inpainting

Motion forecast

$I_{t+1}$

$I_{t-1}$

$I_t$

$t$

VEST-MPI

Quantized scene geometry

Entangled scene representation

Estimated depth

Motion forecast

Disoccluded content inpainting

Novel view

$I_{t+1}$

$I_{t-1}$

$I_t$

$t$

VEST-3D (Ours)

Continuous scene geometry

Disentangled scene representation in 3D space

Entangled scene representation

Novel View

Handling disocclusions

Motion Forecast

VEST

Quantized scene geometry

Multi Plane Images (MPIs)

Generalized MPI Representation

Disentangled scene representation in 3D space

VEST-3D (Ours)

Motion Forecast

Disoccluded Content inpainting

Input Video | Neural Net | **D**ata **A**ssociation | Tracklet Mask (End-to-end Generated) | Tracklet Mask (Linked by **DA**)

plane#1    plane#2

DA

DA

**Efficient Inference**
**Efficient Training**
Inefficient Framework(not end-to-end)
Weak Performance

(a) Frame-level

Clip #1

(Hand-crafted) DA

Clip #2

plane#1    plane#2

**Efficient Inference**
Inefficient Training
Inefficient Framework(partially end-to-end)
**Strong Performance**

(b) Prior Clip-level (VIS Transformer)

Clip #1

RoI-wise

Clip #2

RoI-wise

plane#1    plane#2

**Efficient Inference**
**Efficient Training**
**Efficient Framework(fully end-to-end)**
**Strong Performance**

(c) **EfficientVIS (ours)**